

A note on ADMM

Tri Nguyen
nguyetr9@oregonstate.edu

March 18, 2022

I encounter ADMM several time, did go through the derivation 2 times at least, still I cannot write it down from scratch when a friend of mine asked me about it. This note is to summary my understand about it.

1 Dual Ascend

Firstly, take a look at a common optimization problem with linear equality constraint

$$\underset{x}{\text{minimize}} \quad f(x) \tag{1a}$$

$$\text{subject to} \quad Ax = b, \tag{1b}$$

The first method comes to mind is Lagrangian multiplier. The Lagrangian is

$$L(x, y) = f(x) + y^\top (Ax - b).$$

Optimizing on $L(x, y)$ is easier since it is a unconstrained optimisation problem.

$$\underset{x, y}{\text{minimize}} \quad L(x, y) \tag{2}$$

Solutions (x^*, y^*) of this problem would potentially provide us solution of the original problem.

Block 1.1. *If x^* is a solution of Problem (1), then it would be a part of a stationary point of problem (2), i.e., (x^*, y^*) is a stationary point of problem (2) for some y^* .*

Proof. It might be too much but I will invoke KKT conditions. KKT conditions provides a set necessary conditions for a stationary point of a constrained problem. One of the condition is

$$\frac{\partial L(x^*, y)}{\partial x} = 0$$

Hence, any solution x^* of problem 1 would be stationary point of problem (2). □

Block 1.2 (KKT conditions. (no proof)). *Necessary condition and sufficient condition:*

- *For a constrained optimisation problem with zero duality gap, if x^* and (λ^*, ν^*) are primal and dual optimal then x^*, λ^*, ν^* are satisfied KKT conditions.*
- *For a constrained convex optimisation problem, if x^*, λ^*, ν^* are satisfied KKT conditions, then x^* and (λ^*, ν^*) are primal and dual optimal with zero duality gap.*

Block 1.3 (Strong duality). *These conditions which guarantees strong duality are called constraint qualifications.*

- *If the problem is convex, Slater's condition is a constraint qualification. It requires existence of a strictly feasible point.*

- *Other constraint qualifications.*

So one can work on the Lagrangian problem then pick out the best solution for the original problem. That fact that solving Problem (2) is much easier makes this method a good-to-go method to deal with this constrained problem.

Block 1.4 (Lagrangian multiplier).

2 Derivation

It's not really a derivation but rather a step by step how to get to the procedure above.

1. The augmented Lagrangian function:

$$L_\rho(x, z, y) := f(x) + g(z) + y^T(Ax + Bz - c) + \frac{\rho}{2} \|Ax + Bz - c\|_2^2$$

2. The dual function:

$$g(y) := \inf_{x, z} L_\rho(x, z, y)$$

3. Dual problem:

$$\underset{y}{\text{maximize}} \quad g(y)$$

using gradient ascent, where the gradient respect to y is

$$\nabla g(y) = Ax + Bz - c \quad (\text{prove this})$$

So 3 steps are

$$\begin{aligned} x^{k+1} &\leftarrow \arg \min_x L_\rho(x, z^k, y^k) \\ z^{k+1} &\leftarrow \arg \min_z L_\rho(x^{k+1}, z, y^k) \\ y^{k+1} &\leftarrow y^k + \rho(Ax^{k+1} + Bz^{k+1} - c) \end{aligned}$$

which are equivalent to

$$\begin{aligned} x^{k+1} &\leftarrow \arg \min_x \left(f(x) + \frac{\rho}{2} \|Ax + Bz^k - c + u^k\| \right) \\ z^{k+1} &\leftarrow \arg \min_z \left(g(z) + \frac{\rho}{2} \|Ax^{k+1} + Bz - c + u^k\| \right) \\ u^{k+1} &\leftarrow u^k + Ax^{k+1} + Bz^{k+1} - c \end{aligned}$$

where $u = \frac{1}{\rho}y$.

What if we define residual as $-Ax - Bz + c$?

3 Examples

3.1 Problem 1

$$\begin{aligned} \underset{\mathbf{A}}{\text{minimize}} \quad & \left\| \underline{\mathbf{X}}^{(1)} - (\mathbf{C} \odot \mathbf{B}) \text{diag}(\boldsymbol{\lambda}) \mathbf{A}^T \right\|_{\text{F}}^2 \\ \text{subject to} \quad & \mathbf{A} \geq 0, \mathbf{1}^T \mathbf{A} = \mathbf{1}^T \end{aligned}$$

We can use PGD to solve this, but ADMM is more flexible, and we are not sure which one is faster.

- Step 1: Translate problem to ADMM form

$$\begin{aligned} & \underset{\mathbf{A}, \tilde{\mathbf{A}}}{\text{minimize}} \quad \frac{1}{2} \left\| \underline{\mathbf{X}} - (\mathbf{C} \odot \mathbf{B}) \text{diag}(\boldsymbol{\lambda}) \mathbf{A}^T \right\|_F^2 + I(\tilde{\mathbf{A}} \in \Delta) \\ & \text{subject to} \quad \mathbf{A} = \tilde{\mathbf{A}} \end{aligned}$$

- Step 2: Plug in ADMM updating rule

$$\begin{aligned} \mathbf{A}^{k+1} & \leftarrow \underset{\mathbf{A}}{\text{minimize}} \quad \left(\frac{1}{2} \left\| \underline{\mathbf{X}}^{(1)} - (\mathbf{C} \odot \mathbf{B}) \text{diag}(\boldsymbol{\lambda}) \mathbf{A}^T \right\|_F^2 + \frac{\rho}{2} \left\| \mathbf{A} - \tilde{\mathbf{A}}^k + \mathbf{U}^k \right\|_F^2 \right) \\ \tilde{\mathbf{A}}^{k+1} & \leftarrow \underset{\tilde{\mathbf{A}}}{\text{minimize}} \quad I(\tilde{\mathbf{A}} \in \Delta) + \frac{\rho}{2} \left\| \mathbf{A}^{k+1} - \tilde{\mathbf{A}} + \mathbf{U}^k \right\|_F^2 \\ & = \underset{\tilde{\mathbf{A}}}{\text{minimize}} \quad I(\tilde{\mathbf{A}} \in \Delta) + \frac{\rho}{2} \left\| \tilde{\mathbf{A}} - \mathbf{A}^{k+1} - \mathbf{U}^k \right\|_F^2 \\ & = \mathcal{P}_\Delta(\mathbf{A}^{k+1} + \mathbf{U}^k) \\ \mathbf{U}^{k+1} & \leftarrow \mathbf{U}^k + \mathbf{A}^{k+1} - \tilde{\mathbf{A}}^{k+1} \end{aligned}$$

Solve for \mathbf{A}^T :

$$\begin{aligned} & \nabla = \mathbf{0} \\ \Leftrightarrow & \text{diag}(\boldsymbol{\lambda}) (\mathbf{C} \odot \mathbf{B})^T \left((\mathbf{C} \odot \mathbf{B}) \text{diag}(\boldsymbol{\lambda}) \mathbf{A}^T - \underline{\mathbf{X}}^{(1)} \right) + \rho (\mathbf{A}^T - (\tilde{\mathbf{A}}^k)^T + (\mathbf{U}^k)^T) = \mathbf{0} \\ \Leftrightarrow & \left(\text{diag}(\boldsymbol{\lambda}) (\mathbf{C}^T \mathbf{C} * \mathbf{B}^T \mathbf{B}) \text{diag}(\boldsymbol{\lambda}) + \rho \mathbf{I}_{\text{size}(\mathbf{A}^T, 1)} \right) \mathbf{A}^T = \text{diag}(\boldsymbol{\lambda}) (\mathbf{C} \odot \mathbf{B})^T \underline{\mathbf{X}}^{(1)} + \rho \left((\tilde{\mathbf{A}}^k)^T - (\mathbf{U}^k)^T \right) \\ & \mathbf{G} \mathbf{A}^T = \mathbf{H} \end{aligned}$$

3.2 Example 2

$$\begin{aligned} & \underset{\mathbf{A}}{\text{minimize}} \quad \left\| \underline{\mathbf{X}}^{(1)} - (\mathbf{C} \odot \mathbf{B}) \mathbf{A}^T \right\|_F^2 \\ & \text{subject to} \quad \mathbf{A} \in \mathcal{C} \end{aligned}$$

Step 1: translate to ADMM form. There are several ways to do that, and they are slightly different. Let's stick with the following recipe. The little transpose notation is very easily confusing.

$$\begin{aligned} & \underset{\mathbf{A}, \tilde{\mathbf{A}}}{\text{minimize}} \quad \left\| \underline{\mathbf{X}}^{(1)} - (\mathbf{C} \odot \mathbf{B}) \tilde{\mathbf{A}} \right\|_F^2 + I(\mathbf{A} \in \mathcal{C}) \\ & \text{subject to} \quad \mathbf{A} = \tilde{\mathbf{A}}^T \end{aligned}$$

Step 2:

- Update $\tilde{\mathbf{A}}$,

$$\begin{aligned} \tilde{\mathbf{A}} & \leftarrow \frac{1}{2} \left\| \underline{\mathbf{X}}^{(1)} - (\mathbf{C} \odot \mathbf{B}) \tilde{\mathbf{A}} \right\|_F^2 + \frac{\rho}{2} \left\| \mathbf{A} - \tilde{\mathbf{A}}^T + \mathbf{U} \right\|_F^2 \\ & \leftarrow \frac{1}{2} \left\| \underline{\mathbf{X}}^{(1)} - (\mathbf{C} \odot \mathbf{B}) \tilde{\mathbf{A}} \right\|_F^2 + \frac{\rho}{2} \left\| \tilde{\mathbf{A}} - \mathbf{A}^T - \mathbf{U}^T \right\|_F^2 \\ \nabla_{\tilde{\mathbf{A}}} & = \mathbf{0} \\ \Leftrightarrow & ((\mathbf{C}^T \mathbf{C} * \mathbf{B}^T \mathbf{B}) + \rho \mathbf{I}) \tilde{\mathbf{A}} = (\mathbf{C} \odot \mathbf{B})^T \underline{\mathbf{X}}^{(1)} + \rho (\mathbf{A} + \mathbf{U})^T \\ \Leftrightarrow & \tilde{\mathbf{A}} = [(\mathbf{C}^T \mathbf{C} * \mathbf{B}^T \mathbf{B}) + \rho \mathbf{I}_F]^{-1} [(\mathbf{C} \odot \mathbf{B})^T \underline{\mathbf{X}}^{(1)} + \rho (\mathbf{A} + \mathbf{U})^T] \end{aligned}$$

- Update \mathbf{A} usually with a proximal operator,

$$\mathbf{A} \leftarrow \arg \min_{\mathbf{A}} I(\mathbf{A} \in \mathcal{C}) + \frac{\rho}{2} \left\| \mathbf{A} - \tilde{\mathbf{A}}^T + \mathbf{U} \right\|_F^2$$

- Update dual variable \mathbf{U} ,

$$\mathbf{U} \leftarrow \mathbf{U} + \mathbf{A} - \tilde{\mathbf{A}}^T$$

Does the order of update $\mathbf{A}, \tilde{\mathbf{A}}$ matter? How's about initialization for each sub-problem? Should we init them randomly, or used value from previous iteration?

Small ρ is important to converge to a smaller objective value.